

# Deep Policy

# Dynamic Programming

# for Vehicle Routing Problems



Wouter Kool



Herke van Hoof

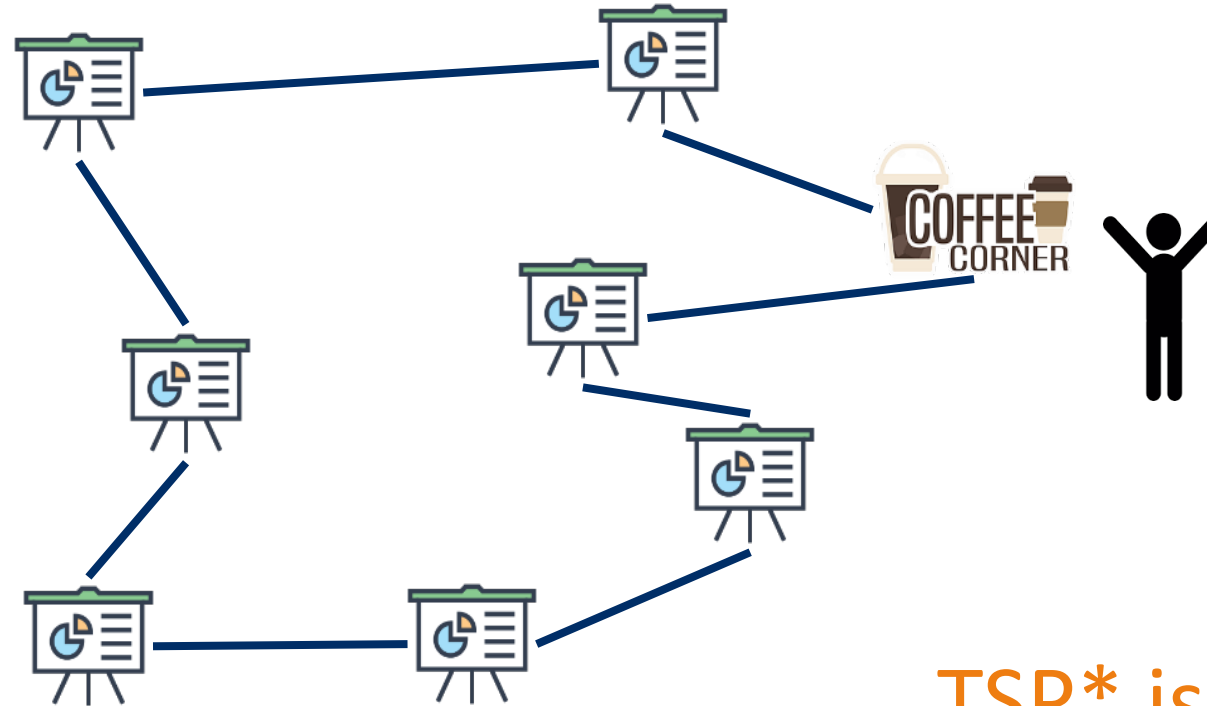


Joaquim Gromicho



Max Welling

# Travelling Scientist Problem (TSP)



TSP\* is (NP-)hard!

Kool et al., 2019

# What does it mean?

Finding *optimal* solutions for *all* problem instances

Finding *acceptable* solutions for *relevant* problem instances

DEEP LEARNING

MISSION:  
~~IMPOSSIBLE~~

\* unless  $P = NP$

MISSION:  
~~IMPOSSIBLE~~

We use HEURISTICS  
can be seen as 'rules of thumb'

'next location should be nearby'

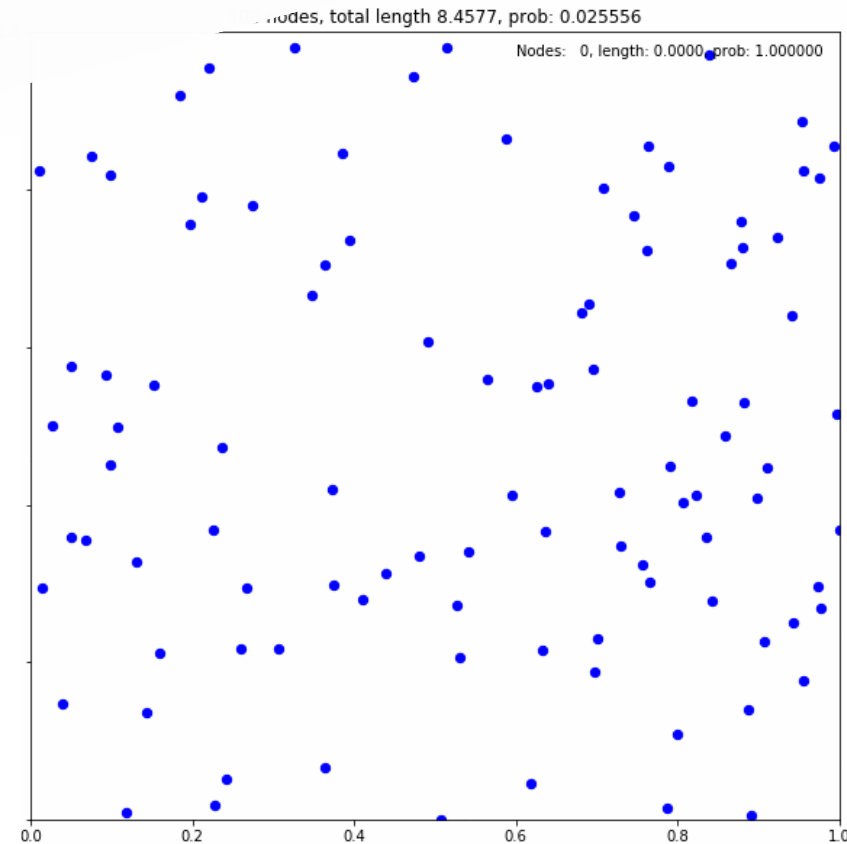
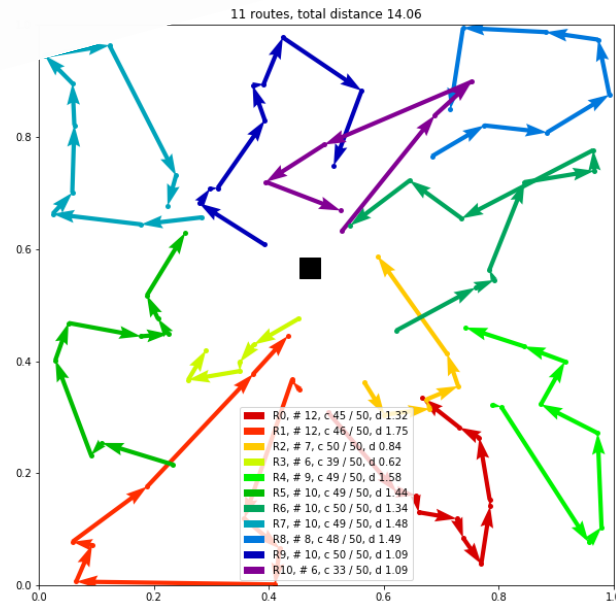
# First try (few years back)

**ATTENTION, LEARN TO SOLVE ROUTING PROBLEMS!**

**Wouter Kool**  
University of Amsterdam  
ORTEC  
w.w.m.kool@uva.nl

**Herke van Hoof**  
University of Amsterdam  
h.c.vanhoof@uva.nl

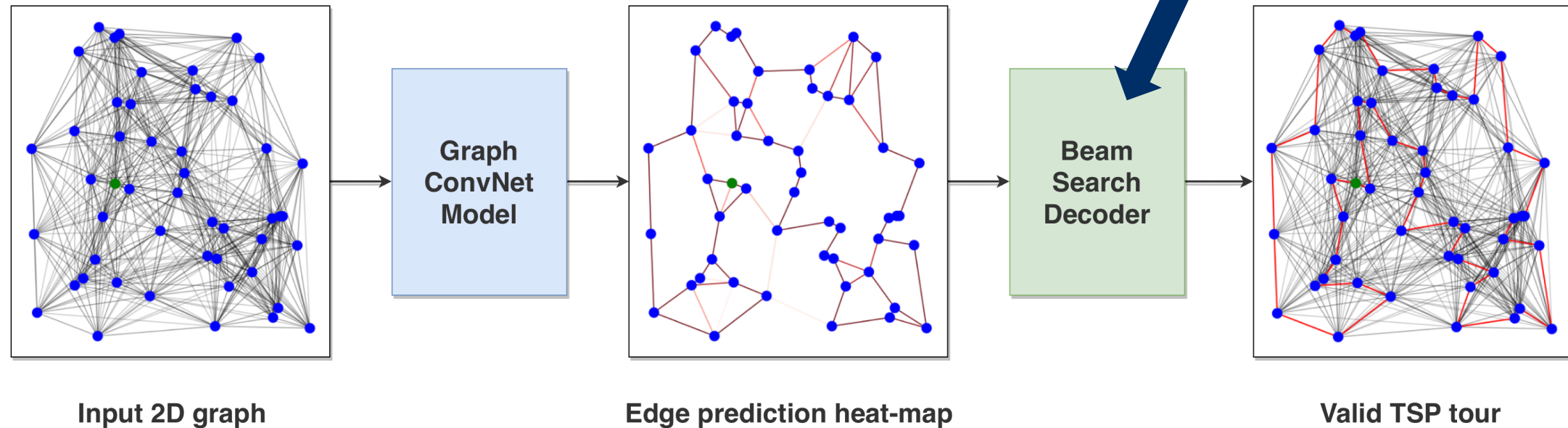
**Max Welling**  
University of Amsterdam  
CIFAR  
m.welling@uva.nl





# Non-autoregressive approach (Joshi et al., 2019)

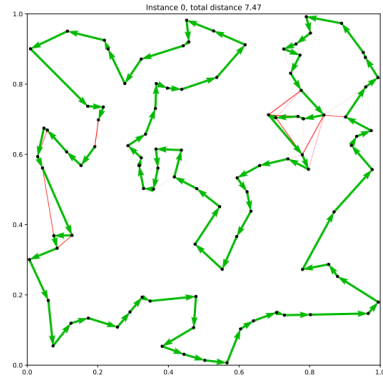
We can do better!



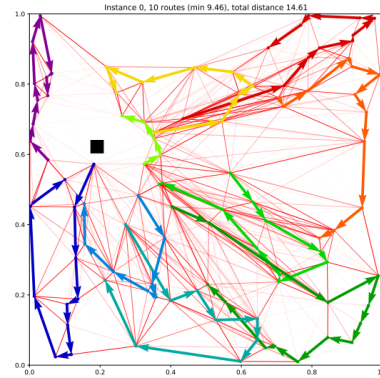
Picture by Joshi et al., 2019

More efficient AND better results!

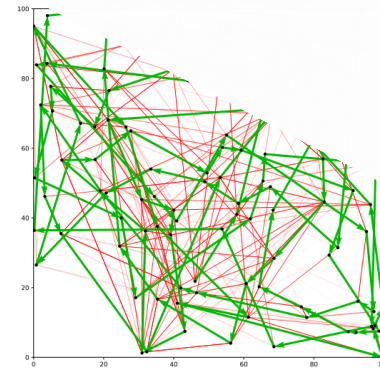
# Bringing us to DPDP



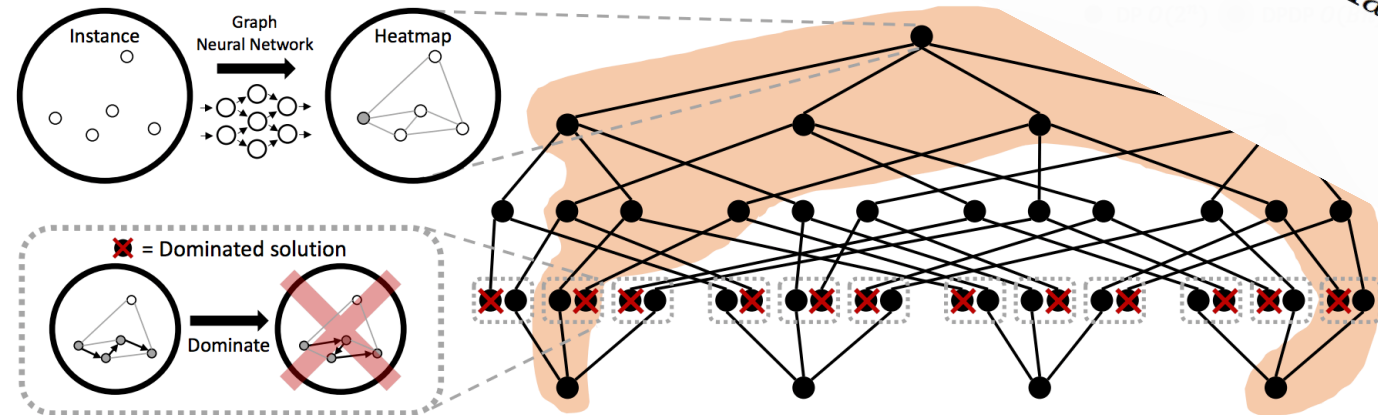
(a) Travelling Salesman Problem



(b) Vehicle Routing Problem



(c) TSP with Time Windows



Deep Policy Dynamic Programming  
for Vehicle Routing Problems  
Wouter Kool<sup>\*1,2</sup>, Herke van Hoof<sup>1</sup>, Joaquim Gromicho<sup>1,2</sup> and Max Welling<sup>1</sup>  
<sup>1</sup>University of Amsterdam  
<sup>2</sup>ORTEC

TL;DR

||

*DPDP is a flexible framework for vehicle routing problems...*

*...that restricts a dynamic program to promising parts of the state space...*

*...using a heatmap of promising edges predicted by a neural network!*

||

# Small side note...

---

If you like 'neural vehicle routing'

# Recent Advances in Deep Learning for Routing Problems

Developing neural network-driven solvers for combinatorial optimization problems such as the Travelling Salesperson Problem have seen a surge of academic interest recently. This blogpost presents a **Neural Combinatorial Optimization** pipeline that unifies several recently proposed model architectures and learning paradigms into one single framework. Through the lens of the pipeline, we analyze recent advances in deep learning for routing problems and provide new directions to stimulate future research towards practical impact.

Chaitanya K. Joshi, Rishabh Anand  
Jan 12, 2022 · 20 min read

<https://www.chaitjo.com/post/deep-learning-for-routing-problems/>

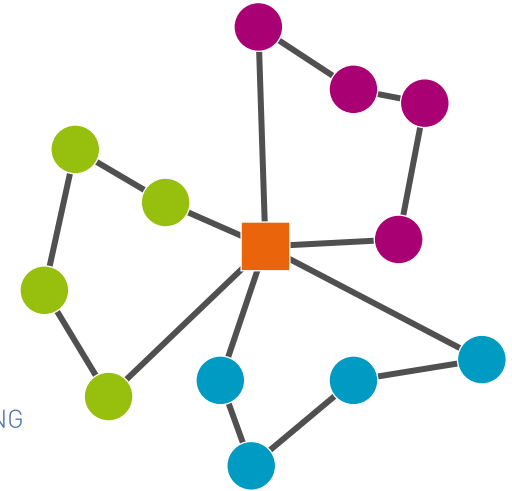
Starting July 1<sup>st</sup>...

Goal: bring together

*Operations Research* and  
*Machine Learning*

to solve a *static* and  
*dynamic* VRP with time  
windows!

# EURO Meets NeurIPS 2022 Vehicle Routing Competition



**ORTEC**

**EAISI** EINDHOVEN  
AI SYSTEMS  
INSTITUTE

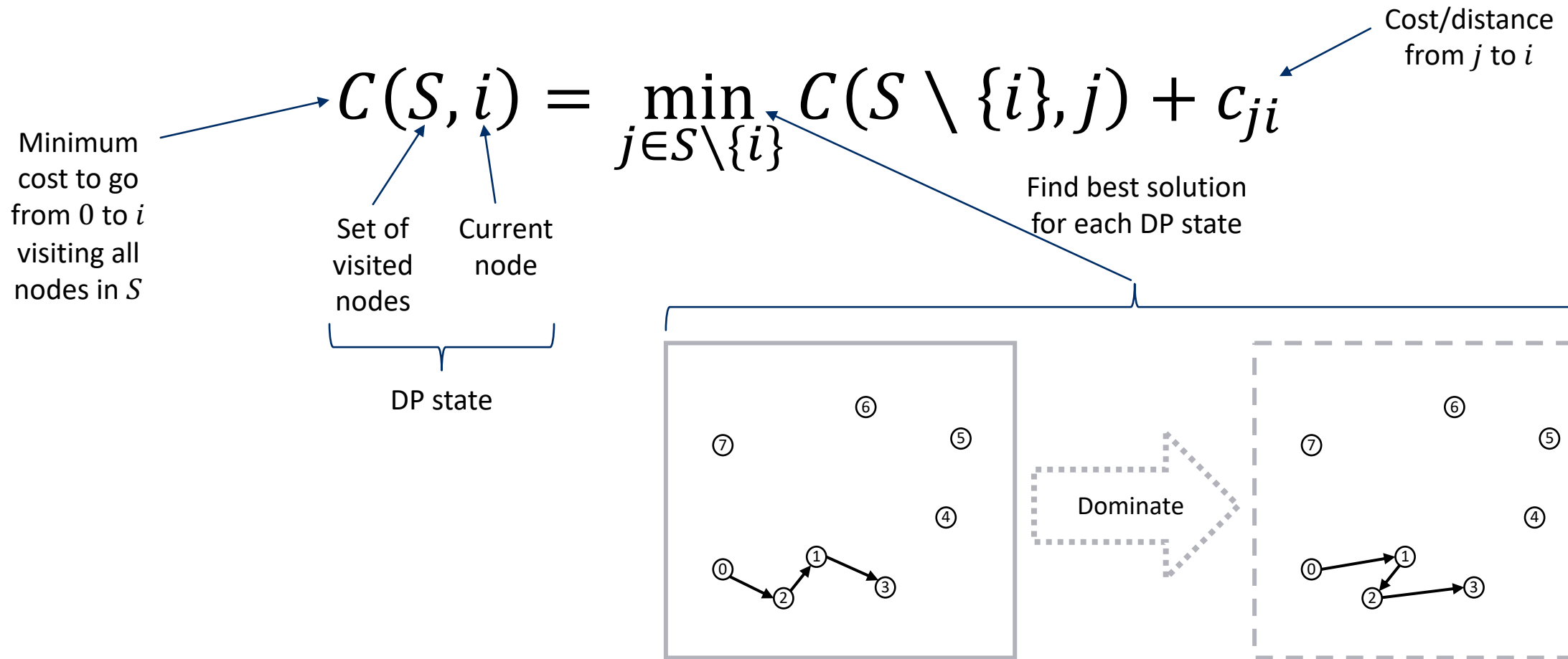
**TU/e**

More info? <https://euro-neurips-vrp-2022.challenges.ortec.com/>

# The technical part

---

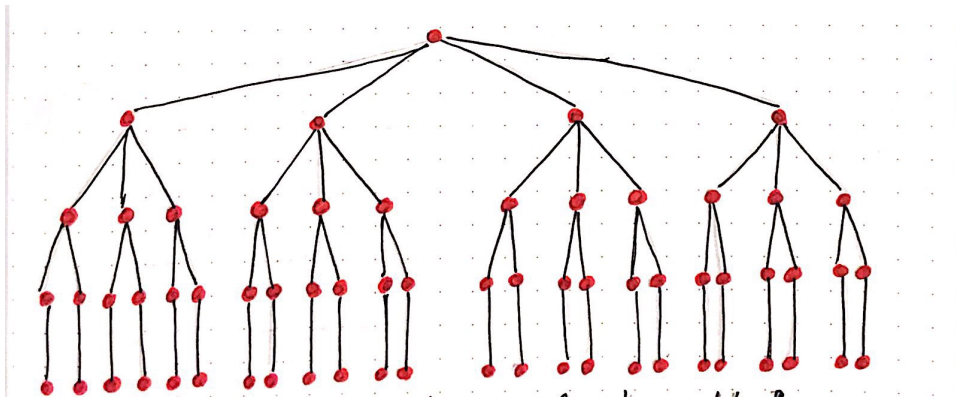
# Held-Karp DP for TSP (Held & Karp, 1962; Bellman, 1962)





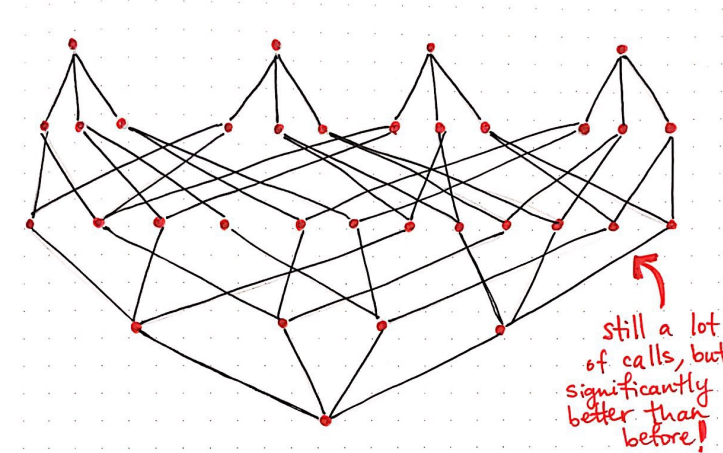
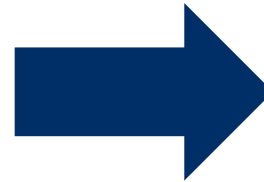
# Held-Karp DP for TSP (Held & Karp, 1962; Bellman, 1962)

Brute-force (forward view)



$O(n!)$  or factorial

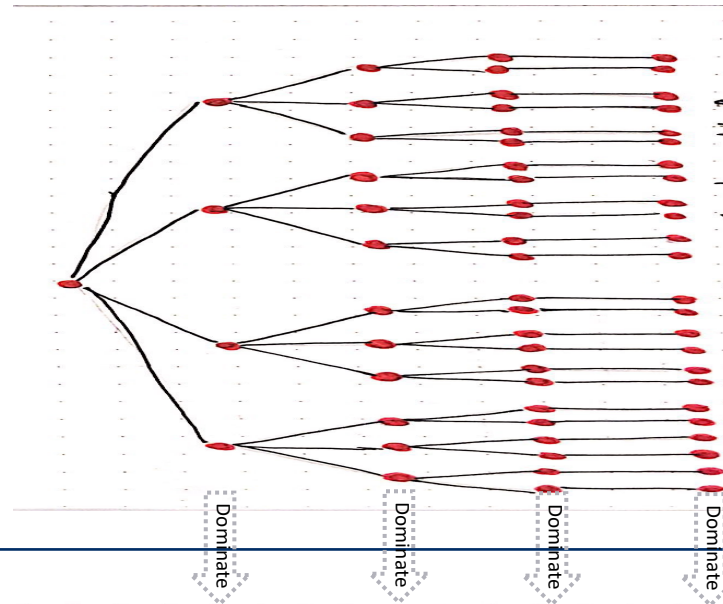
DP (top-down or backward view)



$O(2^n n^2)$  exponential

# Held-Karp DP for TSP (Held & Karp, 1962; Bellman, 1962)

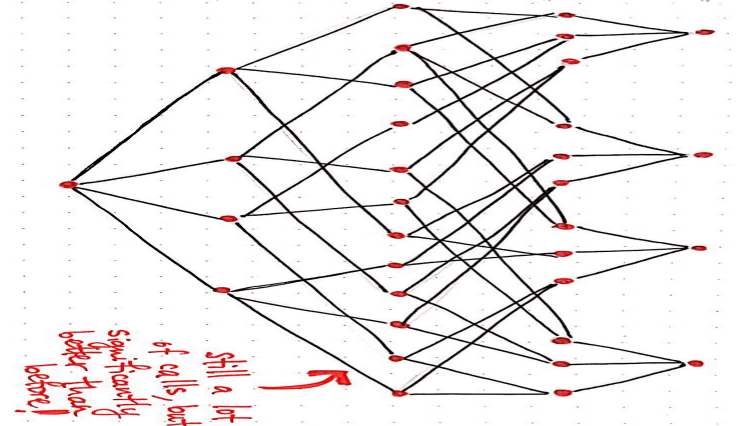
Brute-force  
 $O(n!)$  or factorial



Still impractical!

DP

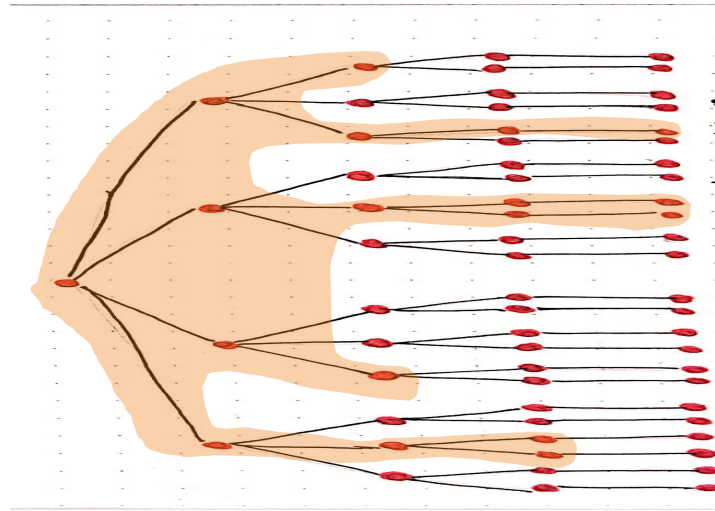
$O(2^n n^2)$  exponential



# Held-Karp DP for TSP (Held & Karp, 1962; Bellman, 1962)

Beam search

$O(Bn)$  or linear

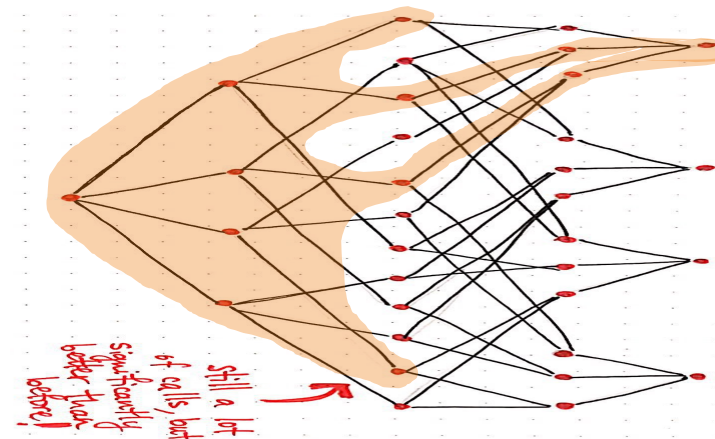


We need a  
*good policy*  
to restrict the  
search space!

Forward view →

Restricted DP

$O(Bn)$  or linear



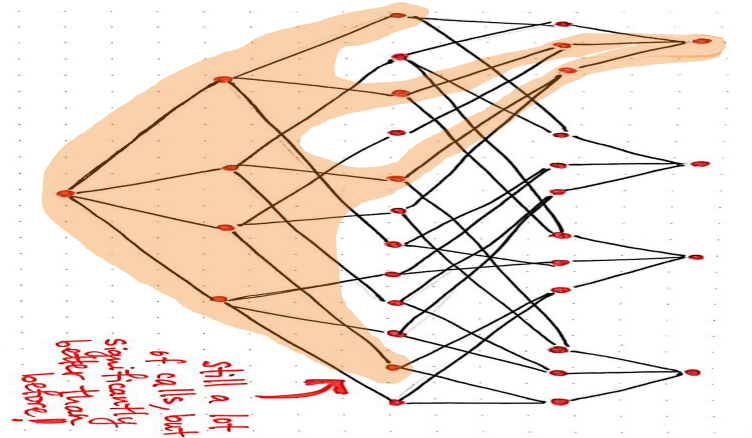
Malandraki & Dial, 1996  
Gromicho et al., 2012

# Deep Policy Dynamic Programming (DPDP)

<https://arxiv.org/abs/2102.11756>

DPDP

$O(Bn)$  or linear



For each iteration

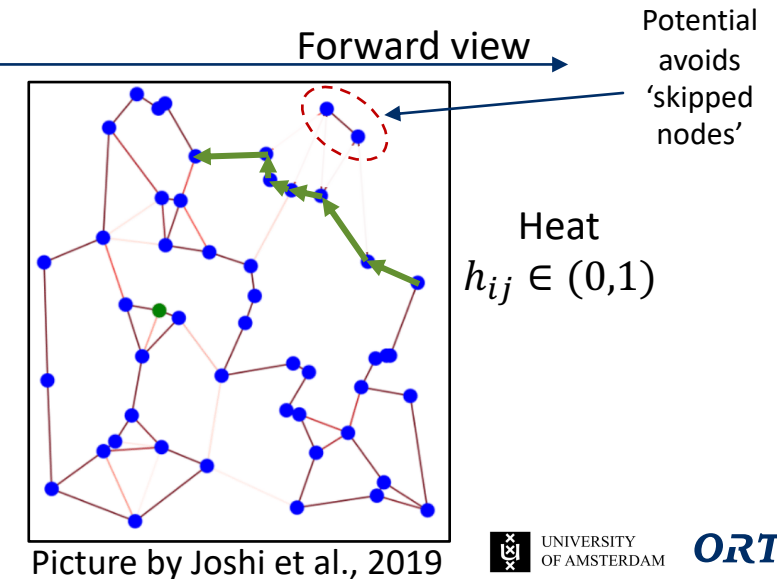
- Expand solutions
- Remove dominated solutions
- Select top  $B$  according to *policy*
- Repeat

*Policy*: select top  $B$  solutions that maximize the score.

$$\text{SCORE} = \text{HEAT} + \text{POTENTIAL}$$

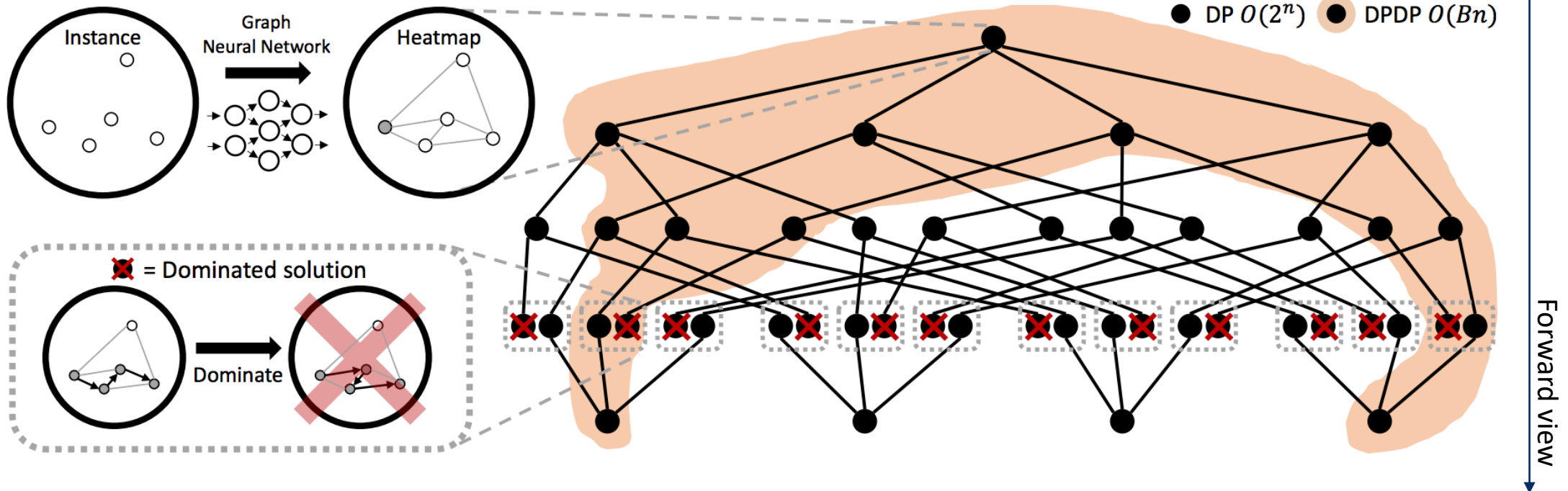
Heat of edges  
in solution

Estimate for  
unvisited nodes  
based on remaining edges



# Deep Policy Dynamic Programming (DPDP)

<https://arxiv.org/abs/2102.11756>



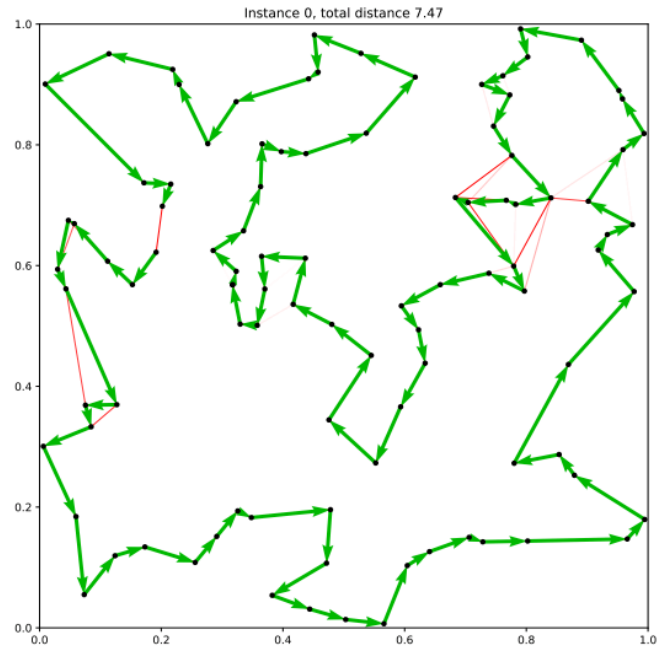
$O(Bn)$  or linear

For each iteration

- Expand solutions
- Remove dominated solutions
- Select top  $B$  according to *policy*
- Repeat



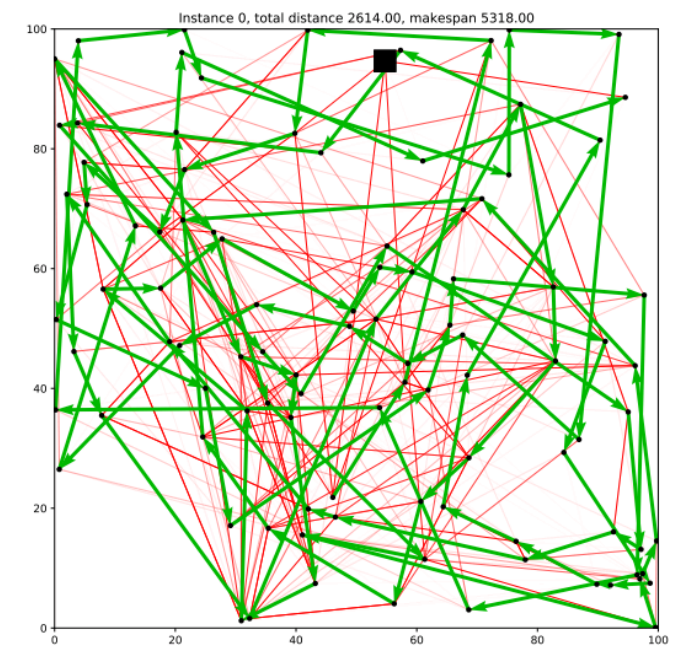
# Experiments



(a) Travelling Salesman Problem



(b) Vehicle Routing Problem



(c) TSP with Time Windows

# Results (TSP/VRP)

Table 1: Mean cost, gap and *total time* to solve 10000 TSP/VRP test instances.

PROBLEM METHOD	TSP100			VRP100		
	COST	GAP	TIME	COST	GAP	TIME
CONCORDE [2]	7.765	0.000 %	6M	15.563	0.000 %	6H11M
HYBRID GENETIC SEARCH [63, 62]						
GUROBI [24]	7.776	0.151 %	31M			
LKH [27]	7.765	0.000 %	42M	15.647	0.536 %	12H57M
GNN HEATMAP + BEAM SEARCH [32]	7.87	1.39 %	40M	16.23	4.28 %	2H
LEARNING 2-OPT HEURISTICS [11]	7.83	0.87 %	41M			
MERGED GNN HEATMAP + MCTS [19]	7.764*	0.04 %	4M + 11M			
ATTENTION MODEL + SAMPLING [36]	7.94	2.26 %	1H			
STEP-WISE ATTENTION MODEL [67]	8.01	3.20 %	29s	16.49	5.96 %	39s
ATTN. MODEL + COLL. POLICIES [34]	7.81	0.54 %	12H	15.98	2.68 %	5H
LEARNING IMPROV. HEURISTICS [66]	7.87	1.42 %	2H	16.03	3.00 %	5H
DUAL-ASPECT COLL. TRANSFORMER [45]	7.77	0.09 %	5H	15.71	0.94 %	9H
ATTENTION MODEL + POMO [37]	7.77	0.14 %	1M	15.76	1.26 %	2M
NEURewriter [9]				16.10	3.45 %	1H
DYNAMIC ATTN. MODEL + 2-OPT [54]				16.27	4.54 %	6H
NEUR. LRG. NEIGHB. SEARCH [30]				15.99	2.74 %	1H
LEARN TO IMPROVE [43]				15.57*	-	4000H
DPDP 10K	7.765	0.009 %	10M + 16M	15.830	1.713 %	10M + 50M
DPDP 100K	7.765	0.004 %	10M + 2H35M	15.694	0.843 %	10M + 5H48M
DPDP 1M				15.627	0.409 %	10M + 48H27M

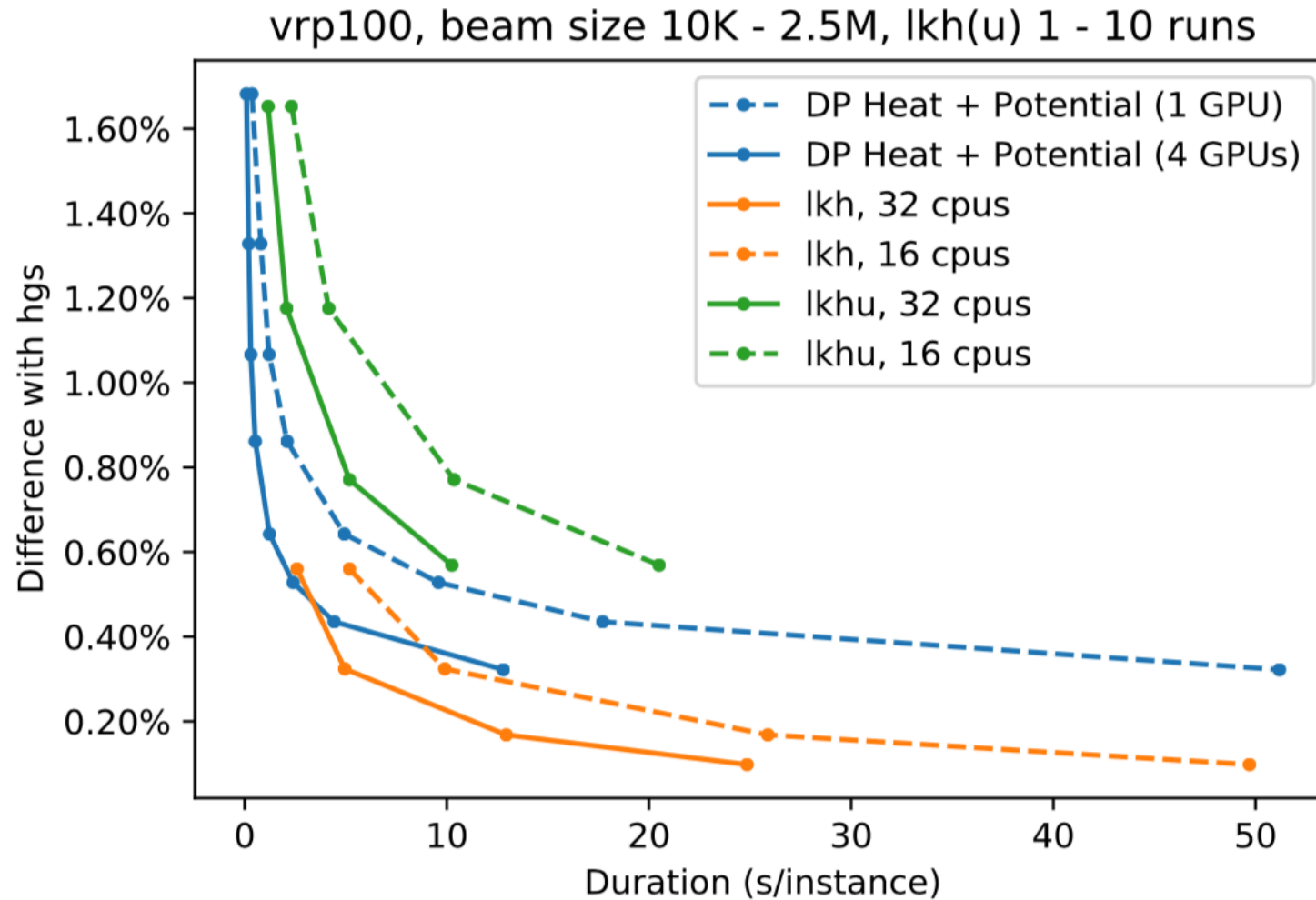
# Results (TSP with time windows)

Table 3: Mean cost, gap and *total time* to solve TSPTW100 instances.

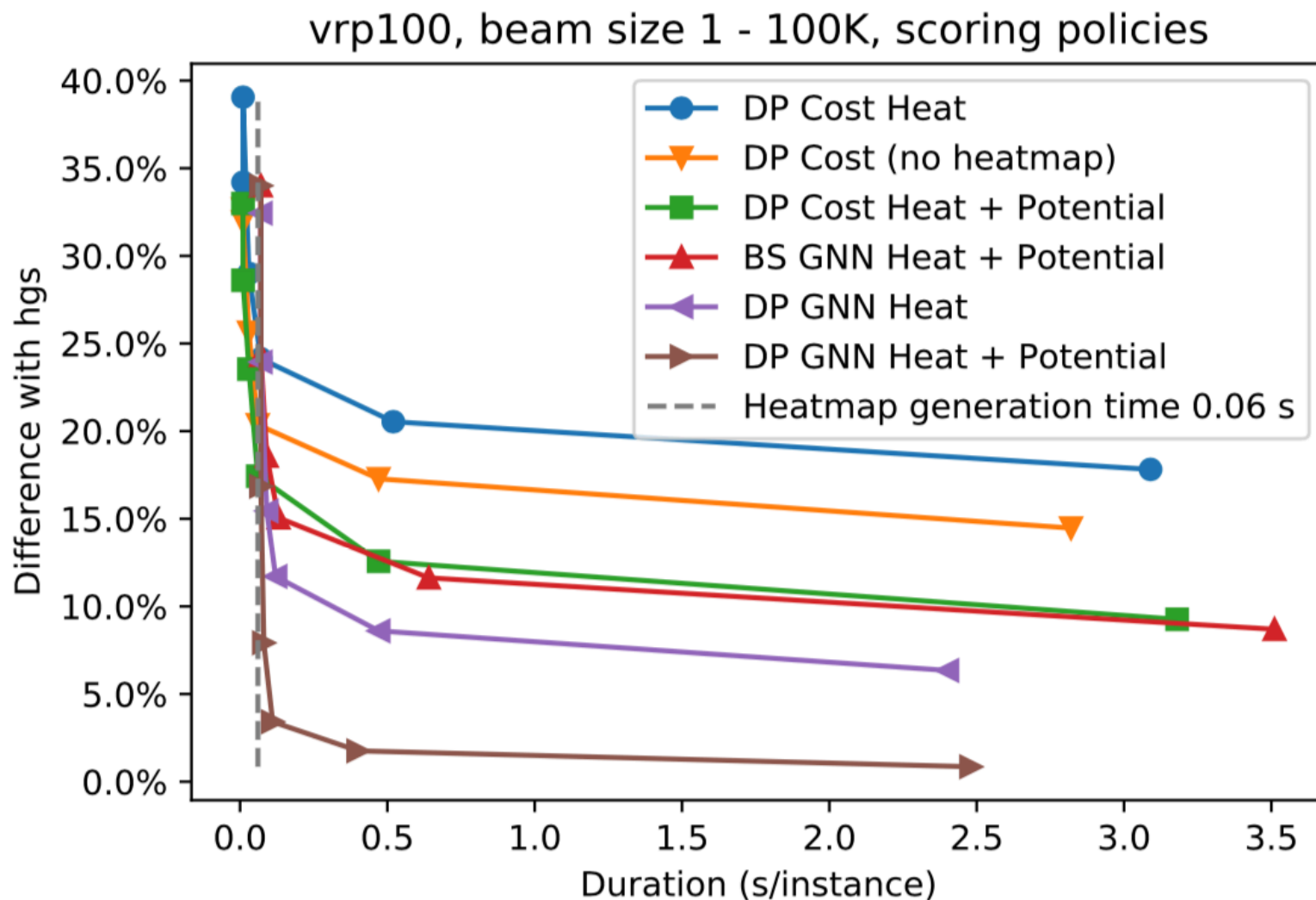
PROBLEM METHOD	SMALL TIME WINDOWS [8] (100 INST.)				LARGE TIME WINDOWS [12] (10K INST.)			
	COST	GAP	FAIL	TIME	COST	GAP	FAIL	TIME
GVNS 30x [12]	5129.58	0.000 %		7s	2432.112	0.000 %		37M15S
GVNS 1x [12]	5129.58	0.000 %		<1s	2457.974	1.063 %		1M4S
LKH 1x [27]	5130.32	0.014 %	1.00 %	5M48S	2431.404	-0.029 %		34H58M
BAB-DQN* [8]	5130.51	0.018 %		25H				
ILDS-DQN* [8]	5130.45	0.017 %		25H				
DPDP 10K	5129.58	0.000 %		6s + 1s	2431.143	-0.040 %		10M + 8M7S
DPDP 100K	5129.58	0.000 %		6s + 1s	2430.880	- 0.051 %		10M + 1h16M



# Quality vs. computation



# Ablations



# Deep Policy Dynamic Programming (DPDP)

<https://arxiv.org/abs/2102.11756>

- DP is flexible framework for many VRP variants e.g. time windows
- Suitable for GPU implementation
- Natural trade-off compute vs. performance -> asymptotically optimal
- Supervised training based on example solutions
- Test time: only evaluate NN once!



